



Published by the Society for Transparency, Openness, and Replication in Kinesiology under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, provided the original author and source are credited.

DOI:

10.51224/cik.2022.49

Subject Areas:

Metascience

Keywords:

Physical therapy, Statistical significance, Inference, Estimation

Author for correspondence:

K. Lohse

lohse@wustl.edu

Editor:

Matthieu Boisgontier



No estimation without inference: A response to the International Society of Physiotherapy Journal Editors

Keith Lohse¹

¹Physical Therapy and Neurology, Washington University
School of Medicine, Saint Louis, MO

The International Society of Physiotherapy Journal Editors (ISPJE) recently published an editorial warning that many of their journals would soon prohibit the use of null hypothesis tests and instead require authors to interpret 95% confidence intervals relative to clinically important values. Although I encourage the reporting of confidence intervals and the discussing of uncertainty in the context of a research question, the ISPJE's proposed ban is illogical and there are several instances of flawed statistical reasoning in the editorial. In brief, the editorial: (1) fails to adequately grapple with the inherent connection between hypothesis testing and estimation, (2) presents several misleading arguments about the perceived flaws of hypothesis tests, and (3) presents an alternative to hypothesis testing that is, in itself, a form of hypothesis test—the minimal effects test—albeit done informally. If the editorials' arguments are taken at face value, then that will lower the statistical literacy in our field and readers will have a flawed understanding of p-values. Further, if the editorials' proposed ban is put into practice, I fear that could decrease the scientific integrity of our research as it removes quantitative benchmarks in favor of a more subjective interpretation of confidence intervals. Ultimately, I think that many of the ISPJE's concerns that led to the editorial are valid, but I think those problems are the result of questionable research practices stemming from poor methodological training for authors, reviewers, and editors. These problems will only be fixed through better and continuing education, not the banning of statistically valid methods.

1. Introduction

Recently, Elkins et al (2022) (hereafter referred to as “the Editorial”) published an editorial on behalf of the International Society of Physiotherapy Journal Editors (ISPJE), recommending that researchers stop using null hypothesis significance tests and adopt “estimation methods”. Further, the editorial warns that this is not merely an idea to consider, but a coming policy of journals: “the [ISPJE] will be expecting manuscripts to use estimation methods *instead* of null hypothesis statistical tests” (emphasis added). However, the Editorial is deeply flawed in its statistical reasoning. If the proposed policies were adopted, they could damage the statistical literacy and scientific integrity of the field.

I detail each of my critiques below, but in short the Editorial: (1) fails to adequately grapple with the inherent connection between hypothesis testing and estimation as methods of statistical inference, (2) presents several misleading arguments about the flaws of statistical significance tests, and (3) presents an alternative that is, in itself, a form of significance testing—the minimal effects test (Murphy & Myors, 1999) (but the alternative does this implicitly and muddles two-sided and one-sided hypothesis testing). Finally, I end with a short list of more urgent problems that the ISPJE could work to address.

I commend the Editorial for encouraging researchers to think deeply about the statistical tools available to them, to consider “practical significance” as well as “statistical significance”, and for bringing important methodological discussions to the forefront of physical therapy research. However, the central argument of the Editorial is illogical and I worry what coming policy changes might mean for how authors interpret their data. I think the antidote to researchers making faulty decisions is not to ban p -values, but to improve education. A rising tide lifts all boats, and if the baseline statistical literacy in our field were higher, authors would make fewer mistakes, reviewers would be more apt to catch remaining mistakes, and readers would be better equipped to make their own conclusions given the available data. Editors then need to hold the line and ensure rigorous review, not ban valid statistical tools.

2. Hypothesis testing and estimation are inescapably intertwined

The Editorial presents hypothesis testing and estimation as two distinct methodological approaches. However, these approaches are two sides of the same coin, as illustrated by a simple example in Figure 1. When a 95% confidence interval excludes the null value, then one can reject the null hypothesis at $p < .05$. This is because hypothesis tests and confidence intervals are based on the same underlying mathematics: e.g., how big is the observed effect relative to the variability we would expect due to sampling? Although typically we think of the null-hypothesis as an assumption of “no effect”, the null hypothesis can assume zero or non-zero effects. So, as shown in Figure 1, we can ascertain the probability of observing the data we did, assuming a null value of 0 or a null value of 1.

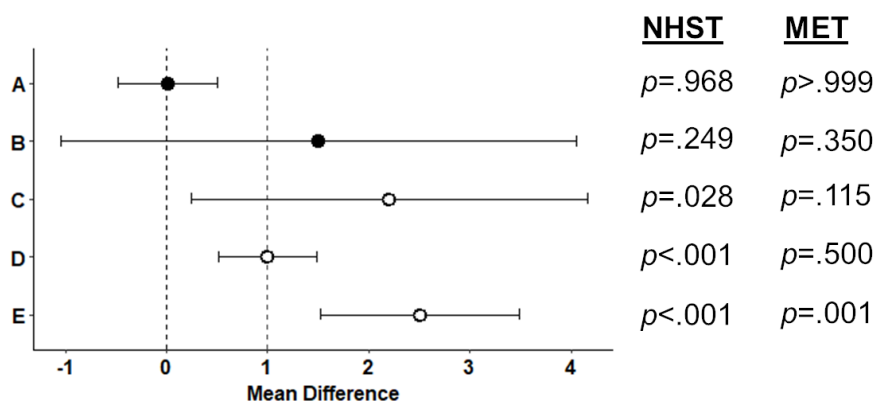


Figure 1: 95% confidence intervals and corresponding p -values for testing $H_0: \Delta = 0$ (NHST, null hypothesis significance testing) and $H_0: \Delta \leq 0$ (MET, a one-sided minimal effects test). Open circles indicate mean differences with $p < 0.05$ for the NHST.

Hypothesis testing and estimation cannot be fully disentangled: estimation (frequentist or Bayesian) asks about *plausible values* of the parameter in the population, hypothesis testing asks about the *plausibility*

of a specific parameter value. These are both inferences, because we are inferring something about the population based on the data in our sample. In the frequentist paradigm, uncertainty in the inference is accounted for with long-run error control; e.g., setting the Type 1 Error rate, $\alpha = 0.05$. We can see this when running simulations as shown in Figure 1A-E: any confidence interval that does not contain zero also has $p < 0.05$ for the null hypothesis significance test (NHST).

The 95% confidence interval shows values in the population that are *compatible* with what we observed in the sample (Rafi & Greenland, 2020). That is, if you move outside of the confidence interval, any of those parameter values (the “true” mean differences; Δ 's) would be statistically different from the mean difference observed in the sample (\bar{X}_d) at the $p < 0.05$ level. Inside of the confidence interval, none of those parameter values would be statistically different ($p > 0.05$) from the observed mean difference. Recall that the p -value is the probability of observing data as extreme or more extreme, assuming that the null hypothesis is true, formally written as $p(\geq \bar{X}_d|H_0)$.

Typically, the null hypothesis significance test (NHST) assumes that the true value in the population is 0 (i.e., $H_0 : \Delta = 0$). The further the sample mean difference is away from 0, the lower the probability of observing that sample mean, if the null hypothesis were true. Importantly, the Editorial does not address the fact that we can set H_0 to be any value. For instance, rather than setting $H_0 : \Delta = 0$ (sometimes referred to as the “nil-hypothesis”) (Cohen, 1994), we can set H_0 equal to any clinically meaningful value of interest. This is referred to as a minimal effects test (or minimum effect test, MET) (Lakens, 2021; Murphy & Myors, 1999). For the sake of argument, let's say this value is 1 in Figure 1. Comparing the confidence intervals to the new null value, you can see that any confidence intervals that only contain values larger than 1 also have a $p < 0.05$ for the minimal effects test (i.e., Figure 1E).¹ Thus, we have both an inference about a specific hypothesis and an estimate in both the NHST and the MET², but the hypothesis test and the estimate are complementary and connected.

3. Misleading arguments about flaws with significance tests

The Editorial bases many arguments on a previous list of perceived problems from Herbert (2019). The Herbert (2019) paper is in itself an editorial that presents informed arguments, but is not an objective demonstration of any mathematical facts. So, reinforcing the Editorial's list through a citation to Herbert (2019) does not provide an evidentiary foundation: it is layering opinion on top of opinion. Second, each of the five “problems” outlined by the Editorial is either not really a problem inherent to p -values or the problem is a true but misleading statement. I address each problem from the Editorial (in quotes) below:

1. “A p -value is not the probability that a hypothesis is (or is not) true.” – This is correct, but it does not follow that this makes p -values, or even statistical significance tests, unhelpful or uninformative. Knowing that the observed data are incompatible with some null value is a crucial step for many research questions. For instance, hypothesis testing in early phase research can help us make decisions about where to direct our resources, starting us down the road of replication and ultimately determining the efficacy and effectiveness of an intervention.

2. A p -value does not constitute evidence – This is an oversimplification and misleading. The Editorial is correct that a single p -value is not strictly speaking “evidence” and cannot tell us about the probability of the null hypothesis being true. However, p -values are still useful tools for making decisions.

Technical definitions of evidence can get a bit complicated and are debated (Goodman & Royall, 1988; Lakens, 2022b; Muff et al., 2022). However, I would invite readers to consider a simple example of absolute probability versus relative probability. If I find that eating green jelly beans reduces post-surgical recovery time for the ACL by 10% relative to controls with $p < 0.05$, then the most likely explanation is still that jelly beans have no effect on recovery and what I observed was chance fluctuation. That is, the null hypothesis is still the most likely explanation even though p was < 0.05 , because the baseline probability of “jelly bean efficacy” is very low and false positives occur 5% of the time when $\alpha = 0.05$. Thus, the p -value is not in itself a measure of evidence, because I would need additional *outside information* in order to change (or not change) my beliefs. As Goodman & Royall (1988) write “The p -value is not adequate for inference because

¹METs are typically directional, using one-sided hypothesis tests (e.g., $H_0 : \leq 1$) whereas NHSTs are often non-directional, using two-sided hypothesis tests (e.g., $H_0 : = 0$). Thus, although the confidence interval for Figure 1A does not contain the null value of 1, the whole of the confidence interval is below 1, thus yielding a non-significant minimal effects test.

²For convenience, I am referring to NHST and MET as separate tests. However, it is more accurate to think of the MET as a type of NHST where you have a one-sided test of a non-zero null value. I use the different terms because readers are likely more familiar with the term NHST when referring to the specific case of $H_0 = 0$ (Cohen, 1994).

the measurement of evidence requires are least three components: the observations, and two competing explanations for how they were produced” (p. 1569; emphasis added).

Some researchers might think of the p -value as evidence against the null specifically, without the need for comparison to a given alternative. But the p -value is calculated assuming that the null is true, so again the Editorial is correct that we cannot simply flip the question around, assume the data, and get the likelihood of the null being true, i.e., $p(\bar{x}_d|H_0) \neq p(H_0|\bar{x}_d)$. To estimate the likelihood of the null hypothesis being true, we would need Bayesian statistics in which we formalize some prior probability about the null hypothesis (Goodman & Royall, 1988). If we have a strong enough prior probability that the null is true, then the current data in the sample may not lead us to change our beliefs based on the posterior distribution, no matter how small the p -value. This was the case in my jelly bean example, where $p < 0.05$ still did not shake my belief in the null hypothesis. For any given prior distribution, however, there is a smaller likelihood of observing highly discrepant effects (e.g., $|\bar{x}_d| \gg 0$), leading to a smaller relative probability of 0 in the posterior distribution compared to the prior distribution.³ Updating the probability of 0 in the posterior distribution reflects rational decision making in daily life. For instance, the first time I find jelly beans reduce recovery time with $p < 0.05$, I might rightly ignore that as a false positive. The fifth time I find jelly beans reduce recovery time with $p < 0.05$, I should take a long hard look at the ingredients and maybe my study procedures; as $p < 0.05$ is not always a sign that the null is wrong, but that some other assumption has been violated.

Still, the p -value does not need to be a measure of evidence for it to be useful. Critically, small p -values are relatively less likely to be observed when the null hypothesis is true compared to when an alternative hypothesis is true. Thus, in a practical sense, a p -values can help us make decisions about what effects to study, assuming that we are testing at least some real effects. As shown in Figure 2A, p -values have a uniform distribution under the null hypothesis, with 5% of p -values necessarily below 0.05. However, if the null is not true, then we will see a shift in the distribution of p -values, with small p -values becoming more common. An example of this is shown in Figure 2B, where the null is false and 34% of p -values are below 0.05. However, correctly rejecting the null hypothesis only 34% of the time is not ideal, so consider Figure 2C, where I have now tripled the sample size and 80% of p -values are below 0.05. That is, with 64 people per group, we now have 80% statistical power to detect a $\Delta = 0.5$.

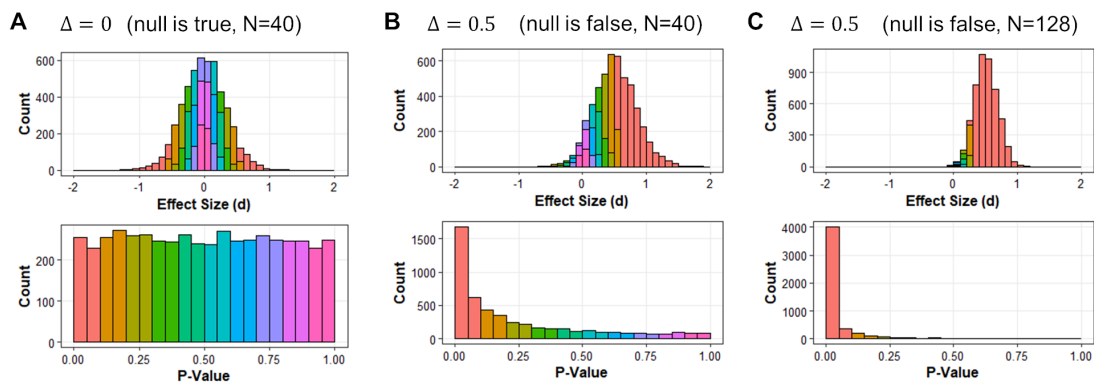


Figure 2: P -values < 0.05 are more likely to occur when the null is false, and critically will only occur 5% of the time when the null is true. Plots show simulated experiments ($k = 5,000$, $\sigma = 1$ for all populations) in which the means of two independent groups are compared using a t -test. In Panel A, the null hypothesis is true and the true difference between population means is 0. In Panel B, the null hypothesis is false and the true difference between population means is 0.5. In Panel C, the null-hypothesis is still false, but I have increased the sample size from 40 to 128, yielding 80% of p -values < 0.05 (i.e., 80% statistical power). Quantiles are color coded with respect to their p -values and effects sizes are given as Cohen’s d .

³For a humorous demonstration see: <https://xkcd.com/1132/>; for a more quantitative visualization of the relationship between priors, p -values, and posteriors see: <https://rpsychologist.com/d3/bayes/>. More technically, the posterior (the updated probability density function after we have seen the data) is proportional to the prior (our expectation before we saw the data) multiplied by the likelihood (which is the probability of the current data given the hypothesis): $posterior \propto likelihood \times prior$.

This is where the concept of a decision is important to distinguish from the term “evidence” (Goodman & Royall, 1988). Without knowing the actual evidence against the null-hypothesis, if I decide to reject the null when $p < 0.05$, then I will only be wrong 5% of the time (i.e., the Type 1 error rate). Similarly, if I have 80% statistical power and a reasonable estimate for the smallest effect size of interest, then I only have a 20% chance of missing an effect of that size (i.e., the Type 2 error rate). Mathematically, these probabilities are robust if we accept the null-hypothesis as true and make minimal other assumptions, which is very helpful when limited outside information is available. See Goodman quoting Neyman and Pearson about hypothesis testing, “Without hoping to know whether each separate hypothesis is true or false, we may search for rules to govern our behaviour with regard to them, in following which we insure that, in the long run of experience, we shall not often be wrong” (Goodman, 1999).

So, p -values are not a measure of evidence, but they are useful tools for helping us make the correct decision. If we want a proper measure of evidence for one hypothesis versus another, then we can do more work, but we also need to make more assumptions and/or bring in outside information. This can be both a feature and bug of *using* hypothesis tests. We can control long run error rates with minimal information, but if we do that so habitually that we forget other information is available, then that is on us not the p -value.

3. “Statistically significant findings are not very replicable.” – This is misleading. First, it is difficult to precisely define replication (Open Science Collaboration, 2015; Patil et al., 2016), but if we think about “being replicable” as the probability that a statistically significant result represents a real, non-zero effect then we would expect more statistically significant findings to “replicate” provided that hypothesis tests have adequate statistical power, researchers have not engaged in p -hacking, there is not selective reporting of results, etc. Thus, not all statistically significant findings will replicate (Scheel et al., 2021), but statistically significant findings in well-designed studies are more likely to replicate (Anderson & Maxwell, 2017; Ioannidis, 2005; Nosek et al., 2022). Second and by any definition, threats to replicability are also going to affect confidence intervals (the Editorial’s proposed solution) as much as they affect p -values, because, again, the p -value is intrinsically linked to the confidence interval. Thus, the Editorial is correct in a practical sense: many statistically significant findings in the current literature do not replicate. However, a lack of replication is the fault of poor study design and questionable research practices, not the use of hypothesis tests as a method of inference.

4. “In most clinical trials, the null hypothesis must be false.” – This is arguably true but very misleading. It is true that real treatment effects are unlikely to be precisely 0 (e.g., they might be +0.001), but it raises the question: Do we really care if the true effect is 0 or 0.001? And will we ever have the statistical precision to discern that difference? All measurement has some error, so I would argue that many effects are functionally 0 even if the (unknowable) true value is not actually zero. But, in a strict mathematical sense I will concede the Editorial is correct, if we accept a hyper-precise definition, the null-hypothesis of $H_0: \Delta = 0.\overline{00}$ will usually be false. However, if we accept that definition, then all point-estimates are false and no value will ever be precisely the minimum clinically important difference either, which is the Editorial’s proposed point-estimate in their alternative.

In response (2022) to an independent critique by Lakens (2022a), this hyper-precise definition does seem to be the argument that the editorial is making.⁴ They claim, “The assertion that the null hypothesis is false in most clinical trials does not require empirical evidence, because it is self-evidently true” and “The null hypothesis may often be approximately true, but it is rarely if ever exactly true”. The Editorial seems to miss the point that the null is a useful model: testing against 0 is still useful for things that are approximately 0. As an analogy, I have successfully found my way many places using maps, but none of those maps was a photo-realistic version of reality.

Scientists are often working on the frontiers of human knowledge; this is costly work where we need to explore a lot of different ideas and many them do not pan out. That is, many tested “effects” are functionally zero (Ioannidis, 2005). So, simply because a point estimate of precisely 0 is unlikely to be true does not mean that it is unhelpful to ask. It should be a very low bar to show that your clinical treatment has a non-zero effect! Further, the Editorial is specifically critiquing this “nil” hypothesis (i.e., $H_0 = 0$), when we could hypothesize any value, or avoid the point-null entirely with a one-sided test (i.e., $H_0 \leq 0$) (Lakens, 2021; Murphy & Myers, 1999). So, if assuming $H_0 = 0$ is not desirable, we can set that

⁴I was very excited to see the Lakens (2022a) commentary and others (Tenan & Caldwell, 2022), and even more excited to see we all largely agree. Interestingly, however, I only became aware of these commentaries after writing my own because I did not see the editorial until it was re-published in Physical Therapy (2022) in June, 2022, whereas my more astute colleagues responded to the original publication in the Journal of Physiotherapy (Elkins et al., 2022), in January 2022. The editorial has been re-published in four different journals to date. While I can appreciate trying to spread one’s message, this creates confusion.

null value to be anything we want (i.e., $H_0 : \Delta \leq 0.4$ m/s for improvement in gait speed, $H_0 : \Delta \leq 30\%$ change on a pain scale, or $H_0 : \Delta \leq 1$ in the hypothetical example in Figure 1).

5. **“Researchers need information about the size of effects.”** – This is a true statement, but it is not a problem with p -values nor null hypothesis significance tests. To my knowledge, no statistician has ever recommended that applied researchers ignore measures of effect size (either raw or standardized). Estimates of effect size are integral to any results section. I would even take this one step further and encourage authors to share their data whenever possible (Borg, Bon, et al., 2020), enabling other researchers to calculate their own effect sizes as there can be limitations with and confusion about standardized effects sizes, and there is no one-size-fits-all solution to effect sizes (Caldwell & Vigotsky, 2020; Levine & Hullett, 2002; McGrath & Meyer, 2006).

4. The Editorial’s “alternative” is a hypothesis test – the Minimal Effects Test

After detailing the potential problems with the NHST, the Editorial proposes an alternative solution in which they encourage authors to compare their 95% confidence interval to some minimum clinically meaningful value (which I will write as δ).⁵ Estimation is a good practice and I would encourage researchers to report 95% confidence intervals and interpret their upper and lower limits in context, when appropriate. However, what the Editorial is suggesting is effectively an MET where $H_0 : \Delta \leq \delta$. That is, if the test is to see if the 95% confidence interval does not contain δ , then that is mathematically equivalent to an MET assuming $H_0 : \Delta \leq \delta$ and finding $p < 0.025$. Note $p < 0.025$, not $p < 0.05$, because most METs are one-sided hypothesis tests whereas confidence intervals are two sided (see Figure 1 and Footnote 1). After heavily critiquing hypothesis testing as a method of inference, the Editorial ends up effectively proposing a hypothesis test. This is clearly an illogical proposition.

I want to emphasize that it is valid for the Editorial to recommend that authors consider their 95% confidence interval relative to some clinically meaningful value. However, this is not an “alternative” to conducting a null hypothesis significance test, it is in fact mathematically identical to conducting a null hypothesis test with a carefully chosen null hypothesis. Both are valid.

I would add, however, that there are also advantages to explicitly framing this as a hypothesis test rather than the informal interpretation of a confidence interval. First, it encourages researchers to explicitly commit to a specific δ while the study is being designed, rather than simply obtaining an estimate of the effect and then comparing it to candidate δ 's post hoc. Second, it requires researchers to think carefully about the direction of the test and the desired α -level, whereas simply invoking a 95% confidence interval implicitly uses a two-tailed test and $\alpha = 0.05$, which may not be best suited to the research question.

Finally, it is also important to stress that history provides us with several examples of how authors will view their data through rose-tinted glasses when quantitative statistical safeguards are removed. For instance, when *Basic and Applied Social Psychology* banned p -values, authors were found to overstate their conclusions well beyond what would have been considered if “statistical significance” had been a benchmark (Fricker Jr. et al., 2019). In sport and exercise science, “magnitude-based inference” was leveraged as a niche method that allowed authors to interpret differences as meaningful when they had very little statistical support (e.g., p 's > 0.25) (Lohse et al., 2020; Sainani, 2018; Sainani et al., 2019). Statistical significance in an NHST does not necessarily need to be the benchmark nor 0.05 the default value (Amrhein & Greenland, 2018; Benjamin et al., 2018; Lakens et al., 2018; McShane et al., 2019), but it is always important to have a statistically sound framework for dealing with uncertainty.

5. Virtues of hypothesis testing

One of the great virtues of null hypothesis significance testing is Type I error control while making minimal assumptions about the nature of the data or the world at large. If we set $\alpha = 0.05$, then we can be confident we will only get data greater than or equal to what we observed 5% of the time when the null is true. Importantly, this works for a wide range of statistics and types of tests, including F - and χ^2 -statistics that have multiple degrees of freedom from models asking questions about multiple effects simultaneously. For instance, in a randomized controlled trial with three arms, I could conduct an omnibus F -test and obtain a p -value to see if there is any evidence of a difference between groups overall, before conducting

⁵I caution that it is difficult to find a single measure of δ ; it changes as a function of the study population, the study context, and has its own uncertainty due to sampling error (Dabija & Jain, 2019; Tenan & Caldwell, 2022).

additional post-hoc tests to compare specific groups. This situation is not covered by the Editorial and the Editorial's confidence interval alternative is not easily applied here, although one could plausibly adjust the width of the confidence intervals to control for multiple comparisons.

6. Bigger threats to statistical integrity

Misinterpretation and misuse of p -values are threats to statistical integrity. However, questionable research practices such as p -hacking, sub-group analyses, flexible stopping rules, selective exclusion of outliers, selective reporting, and hypothesizing after results are known (HARKing) are much larger threats (Kerr, 1998; Rosenthal, 1979; Simmons et al., 2011, 2013; Sun et al., 2012). Furthermore, these questionable research practices have consistently negative consequences regardless of the method of inference. For instance, although the term " p -hacking" connotes the NHST, these questionable research practices pose an equal threat to confidence intervals because again confidence intervals and p -values are based on the same underlying mathematics. Similarly, switching to a fully Bayesian method of analysis is not an antidote for poor study design, small samples, and questionable research practices. As others have argued (Borg, Lohse, et al., 2020; Leek & Peng, 2015), p -values get a disproportionate amount of attention in popular discussions of research methodology. I encourage the ISPJE to instead focus their attention on methods for improving data/code sharing, transparency, and replicability through tools like preregistration, results-blind review, registered reports, or even "data papers" whose primary function is to report a study and archive the data, without drawing inferences from limited samples.

It is entirely valid to say that p -values are often mis-used and mis-interpreted, and "statistical significance" may not ultimately be the best term for applied researchers to use (Wasserstein et al., 2019). However, it is incorrect to present these human errors as inherent flaws in hypothesis testing. For instance, if someone mis-interprets $p > 0.05$ as evidence of "no difference", then I would argue the correct action is to teach them about equivalence tests and non-inferiority designs, not ban p -values. Similarly, there are times when Bayesian inference is what authors are really interested in (e.g., what is the probability that the null is true, given the evidence?), and in those cases Bayesian inference can and should be used. However, Bayesian analysis is not a panacea and needs to be used thoughtfully like any statistical tool. So, although a simple heuristic of $p < 0.05$ may well be overused as "the" test in physical therapy research, frequentist hypothesis tests are still valid and useful tools for physical therapy researchers. Moreover, the scientific integrity of the field has much larger concerns, and both p -values and confidence intervals will be corrupted by p -hacking, under-powered subgroup analyses, surrogate outcomes, and other questionable research practices.

In conclusion, I agree with the Editorial on the importance of reporting effect sizes and interpreting them in context. However, the Editorial makes numerous statistical *faux pas* that could harm the statistical literacy in our field, if readers take them at face value, and harm the scientific integrity of our field, if put into editorial practice.

7. Additional Information

(a) Data Accessibility

R code for all analyses and simulations presented in this commentary are included as a digital supplement on SportRxiv (<https://sportrxiv.org/index.php/server/preprint/view/178/version/211>).

(b) Conflict of Interest

Author has no conflicts of interest to declare.

(c) Funding

None.

(d) Acknowledgments

I would like to thank Dr. Emma Johnson, Dr. Kristin Sainani, and two anonymous reviewers for their detailed comments on earlier drafts of this commentary.

(e) Preprint

The pre-publication version of this manuscript can be found on SportRxiv (DOI: <https://doi.org/10.51224/SRXIV.178>).

8. References

- Amrhein, V., & Greenland, S. (2018). Remove, rather than redefine, statistical significance. *Nature Human Behavior*, 2, 4–4. <https://doi.org/10.1038/s41562-017-0224-0>
- Anderson, S. F., & Maxwell, S. E. (2017). Addressing the “replication crisis”: Using original studies to design replication studies with appropriate statistical power. *Multivariate Behavioral Research*, 52, 305–324. <https://doi.org/10.1080/00273171.2017.1289361>
- Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E. J., Berk, R., ..., & Johnson, V. E. (2018). Redefine statistical significance. *Nature Human Behavior*, 2, 6–10. <https://doi.org/10.1038/s41562-017-0189-z>
- Borg, D. N., Bon, J., Sainani, K. L., Baguley, B. J., Tierney, N., & Drovandi, C. (2020). Sharing data and code: A comment on the call for the adoption of more transparent research practices in sport and exercise science. *SportRxiv*. <https://doi.org/10.31236/osf.io/ftdgj>
- Borg, D. N., Lohse, K. R., & Sainani, K. L. (2020). Ten common statistical errors from all phases of research, and their fixes. *PM&R*, 12, 610–614. <https://doi.org/10.1002/pmrj.12395>
- Caldwell, A., & Vigotsky, A. D. (2020). A case against default effect sizes in sport and exercise science. *PeerJ*, 8, e10314. <https://doi.org/10.7717/peerj.10314>
- Cohen, J. (1994). The earth is round ($p < .05$). *American Psychologist*, 49, 997–1003. <https://doi.org/10.1037/0003-066X.49.12.997>
- Dabija, D. I., & Jain, N. B. (2019). Minimal clinically important difference of shoulder outcome measures and diagnoses: A systematic review. *American Journal of Physical Medicine & Rehabilitation*, 98, 671–676. <https://doi.org/10.1097/phm.0000000000001169>
- Elkins, M. R., Pinto, R. Z., Verhagen, A., Grygorowicz, M., Soderlund, A., Guemann, M., Gomez-Conesa, A., Blanton, S., Brismée, J. M., Agarwal, S., Jette, A., Karstens, S., Harms, M., Verheyden, G., & Sheikh, U. (2022). Statistical inference through estimation: Recommendations from the International Society of Physiotherapy Journal Editors. *Physical Therapy*, 102, pzac066. <https://doi.org/10.1016/j.jphys.2021.12.001>
- Elkins, M. R., Pinto, R. Z., Verhagen, A., Grygorowicz, M., Söderlund, A., Guemann, M., Gómez-Conesa, A., Blanton, S., Brismée, J. M., Agarwal, S., Jette, A., Harms, M., Verheyden, G., & Sheikh, U. (2022). Correspondence: Response to Lakens. *Journal of Physiotherapy*, 68, 214. <https://doi.org/10.1016/j.jphys.2022.06.003>
- Elkins, M. R., Pinto, R. Z., Verhagen, A., Grygorowicz, M., Söderlund, A., Guemann, M., Gómez-Conesa, A., Blanton, S., Brismée, J. M., Ardern, C., Agarwal, S., Jette, A., Karstens, S., Harms, M., Verheyden, G., & Sheikh, U. (2022). Statistical inference through estimation: Recommendations from the International Society of Physiotherapy Journal Editors. *Journal of Physiotherapy*, 68(1), 1–4. <https://doi.org/10.1016/j.jphys.2021.12.001>
- Fricker Jr., R. D., Burke, K., Han, X., & Woodall, W. H. (2019). Assessing the statistical analyses used in Basic and Applied Social Psychology after their p-value ban. *The American Statistician*, 73, 374–384. <https://doi.org/10.1080/00031305.2018.1537892>
- Goodman, S. N. (1999). Toward evidence-based medical statistics. 1: The P value fallacy. *Annals of Internal Medicine*, 130, 995–1004. <https://doi.org/10.7326/0003-4819-130-12-199906150-00008>
- Goodman, S. N., & Royall, R. (1988). Evidence and scientific research. *American Journal of Public Health*, 78, 1568–1574. <https://doi.org/10.2105/ajph.78.12.1568>
- Herbert, R. (2019). Research note: Significance testing and hypothesis testing: Meaningless, misleading and mostly unnecessary. *Journal of Physiotherapy*, 65, 178–181. <https://doi.org/10.1016/j.jphys.2019.05.001>
- Ioannidis, J. P. (2005). Why most published research findings are false. *PLoS Medicine*, 2, e124. <https://doi.org/10.1371/journal.pmed.0020124>
- Kerr, N. L. (1998). HARKing: Hypothesizing after the results are known. *Personality and Social Psychology Review*, 2, 196–217. https://doi.org/10.1207/s15327957pspr0203_4
- Lakens, D. (2021). The practical alternative to the p value is the correctly used p value. *Perspectives on Psychological Science*, 16, 639–648. <https://doi.org/10.1177/1745691620958012>
- Lakens, D. (2022a). Correspondence: Reward, but do not yet require, interval hypothesis tests. *Journal of Physiotherapy*, 68, 213–214. <https://doi.org/10.1016/j.jphys.2022.06.004>
- Lakens, D. (2022b). Why P values are not measures of evidence. *Trends in Ecology & Evolution*, 37, 289–290. <https://doi.org/10.1016/j.tree.2021.12.006>

- Lakens, D., Adolphi, F. G., Albers, C. J., Anvari, F., Apps, M. A., Argamon, S. E., ..., & Zwaan, R. A. (2018). Justify your alpha. *Nature Human Behavior*, 2, 168–171. <https://doi.org/10.1038/s41562-018-0311-x>
- Leek, J. T., & Peng, R. D. (2015). Statistics: P values are just the tip of the iceberg. *Nature*, 520, 612–612. <https://doi.org/10.1038/520612a>
- Levine, T. R., & Hullett, C. R. (2002). Eta squared, partial eta squared, and misreporting of effect size in communication research. *Human Communication Research*, 28, 612–625. <https://doi.org/10.1111/j.1468-2958.2002.tb00828.x>
- Lohse, K. R., Sainani, K. L., Taylor, A., Butson, M. L., Knight, E. J., & Vickers, A. J. (2020). Systematic review of the use of “magnitude-based inference” in sports science and medicine. *PLoS One*, 15, 0235318. <https://doi.org/10.1371/journal.pone.0235318>
- McGrath, R. E., & Meyer, G. J. (2006). When effect sizes disagree: the case of r and d. *Psychological Methods*, 11, 386. <https://doi.org/10.1037/1082-989x.11.4.386>
- McShane, B. B., Gal, D., Gelman, A., Robert, C., & Tackett, J. L. (2019). Abandon statistical significance. *The American Statistician*, 73, 235–245. <https://doi.org/10.1080/00031305.2018.1527253>
- Muff, S., Nilsen, E. B., O’Hara, R. B., & Nater, C. R. (2022). Response to “Why P values are not measures of evidence” by D. Lakens. *Trends in Ecology & Evolution*, 37, 291–292. <https://doi.org/10.1016/j.tree.2022.01.001>
- Murphy, K. R., & Myers, B. (1999). Testing the hypothesis that treatments have negligible effects: Minimum-effect tests in the general linear model. *Journal of Applied Psychology*, 84, 234–248. <https://doi.org/10.1037/0021-9010.84.2.234>
- Nosek, B. A., Hardwicke, T. E., Moshontz, H., Allard, A., Corker, K. S., Dreber, A., Fidler, F., Hilgard, J., Struhl, M. K., Nuijten, M. B., Rohrer, J. M., Romero, R., Scheel, A. M., Scherer, L. D., Schönbrodt, B. A., Nosek, F. D., & Vazire, S. (2022). Replicability, robustness, and reproducibility in psychological science. *Annual Review of Psychology*, 73, 719–748. <https://doi.org/10.1146/annurev-psych-020821-114157>
- Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science*, 349, 4716. <https://doi.org/10.1126/science.aac4716>
- Patil, P., Peng, R. D., & Leek, J. T. (2016). What should researchers expect when they replicate studies? A statistical view of replicability in psychological science. *Perspectives on Psychological Science*, 11, 539–544. <https://doi.org/10.1177/1745691616646366>
- Rafi, Z., & Greenland, S. (2020). Semantic and cognitive tools to aid statistical science: replace confidence and significance by compatibility and surprise. *BMC Medical Research Methodology*, 20, 244. <https://doi.org/10.1186/s12874-020-01105-9>
- Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological Bulletin*, 86. <https://doi.org/10.1037/0033-2909.86.3.638>
- Sainani, K. L. (2018). The Problem with “Magnitude-based Inference.” *Medicine and Science in Sports and Exercise*, 50, 2166–2176. <https://doi.org/10.1249/mss.0000000000001645>
- Sainani, K. L., Lohse, K. R., Jones, P. R., & Vickers, A. (2019). Magnitude-based inference is not Bayesian and is not a valid method of inference. *Scandinavian Journal of Medicine & Science in Sports*, 29, 1428. <https://doi.org/10.1111/sms.13491>
- Scheel, A. M., Schijen, M. R. M. J., & Lakens, D. (2021). An excess of positive results: Comparing the standard Psychology literature with Registered Reports. *Advances in Methods and Practices in Psychological Science*, 4, 25152459211007468. <https://doi.org/10.1177/25152459211007467>
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22, 11, 1359–1366. <https://doi.org/10.1177/0956797611417632>
- Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2013). Life after p-hacking. *Meeting of the Society for Personality and Social Psychology*, 17–19. <https://dx.doi.org/10.2139/ssrn.2205186>
- Sun, X., Briel, M., Busse, J. W., You, J. J., Akl, E. A., Mejza, F., ..., & Guyatt, G. H. (2012). Credibility of claims of subgroup effects in randomised controlled trials: Systematic review. *BMJ*, 344. <https://doi.org/10.1136/bmj.e1553>
- Tenan, M. S., & Caldwell, A. R. (2022). Confidence intervals and smallest worthwhile change are not a panacea: A response to the International Society of Physiotherapy Journal Editors. *Communications in Kinesiology*, 1, 4, 1–10. <https://doi.org/10.51224/cik.2022.45>
- Wasserstein, R. L., Schirm, A. L., & Lazar, N. A. (2019). Moving to a world beyond “p < 0.05.” *The American Statistician*, 73, 1–19. <https://doi.org/10.1080/00031305.2019.1583913>